

**DM566/DM868/DM870: Data Mining and Machine Learning**  
Spring term 2019

**Exercise 12: Hierarchical Clustering, Outlier Detection**

**Exercise 12-1 OPTICS Plot**

(a) For the data below we got computed the reachability diagram to the right.



With a naïve understanding of hierarchical clustering, wouldn't we have expected three valleys in the plot? Explain, why this is not the case and why the plot, instead, looks as it does and accurately describes the density structure of the data.

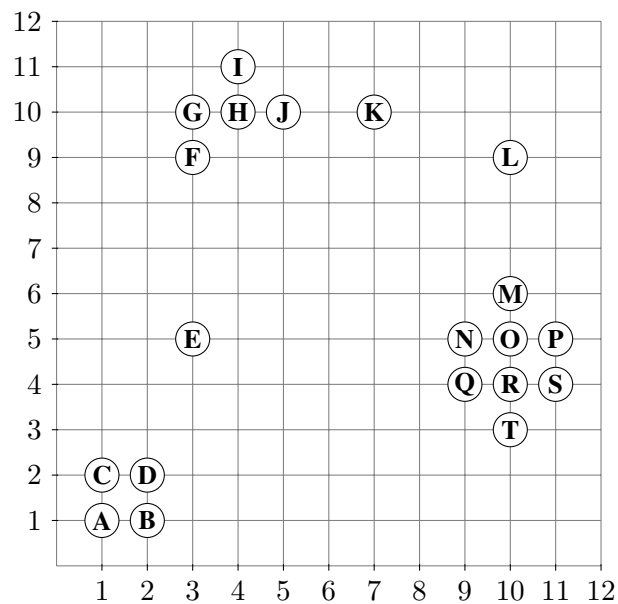
(b) For this dataset (left) we have the reachability plot (right).



Mark in the reachability plot which areas relate to the clusters *A*, *B*, *C*, *D*, and *E*.

### Exercise 12-2      Outlier Scores

Given the following 2 dimensional data set:



As distance function, use Manhattan distance  $L_1(a, b) := |a_1 - b_1| + |a_2 - b_2|$ .

Compute the following (without including the query point when determining the  $k$ NN):

- LOF using  $k = 2$  for the points  $E$ ,  $K$  and  $O$ .
- LOF using  $k = 4$  for the points  $E$ ,  $K$  and  $O$ .
- $k$ NN distance using  $k = 2$  for all points.
- $k$ NN distance using  $k = 4$  for all points.
- aggregated  $k$ NN distances for  $k = 2$  and  $k = 2$  for all points  
(aggregated  $k$ NN distance = sum of the distances to all the  $k$ NN!)

**Exercise 12-3      Evaluation of Outlier Scores**

A data set with known outliers + was evaluated using two outlier detection methods  $S_1$  and  $S_2$ . The results of the methods are given in the table below:

Object	$A_1$	$A_2$	$A_3$	$A_4$	$A_5$	$A_6$	$A_7$	$A_8$	$A_9$	$A_{10}$
Label	–	–	–	–	+	–	–	–	+	–
$S_1$	1.0	1.1	1.1	1.3	3.0	2.0	1.5	0.9	1.4	1.2
$S_2$	.80	.80	.10	.81	.89	.50	.50	.91	.90	.20

Evaluate both outlier detection methods  $S_1$  and  $S_2$  using the following metrics:

- Precision, Recall and F-Measure, assuming that the top  $k = 2$  ranked outliers were classified as outliers.
- Average Precision for  $k = 1 \dots 4$ , assuming that the top  $k$  ranked outliers were classified as outliers.
- Draw the ROC curve, and compute the area under curve (AUC) measure.

**Exercise 12-4      Outlier Detection – Practical**

Bring your laptop for this interactive lab session in the exercise class.

- Work with some toolbox for data exploration (e.g., R, Python, ELKI) to try different outlier detection algorithms (e.g., knn, LOF) on some dataset (e.g., “3 clusters and noise 2d” from e-learn).
- How does the behavior change with the choice of the neighborhood size?
- Run OPTICS on the same dataset. Imagine, you would not know how the dataset looks like. What could you learn about the clusters and outliers in the dataset?